# INFORMATION RETRIEVAL

Week 4 – Tolerant Retrieval

21.03.2025

l

Severin Mills

### Today



- Discussion
- Questions

- B+-Trees
- Wildcard Queries
- Permuterm Index
- Spelling correction

Exercise 3: Tolerant • Retrieval



4

### Kahoot / Exam questions

## Discussion

	Wahr	Falsch	
	0	$\bigcirc$	Stemming should be invoked at indexing time processing a query.
	$\bigcirc$	$\bigcirc$	In a Boolean retrieval system, stemming never
×	$\bigcirc$	$\bigcirc$	Stemming increases the size of the vocabulary
	0	$\bigcirc$	In a Boolean retrieval system, stemming never precision.
03 202	5 –		



### Discussion

- Which of the following bit strings contain only valid UTF-8 encoded data?
  - 01001000

Yes it does, 8-bit string starting with 0

### Discussion

- Which of the following bit strings contain only valid UTF-8 encoded data?
  - 11001110 10110001

Yes it does, two 8-bit strings, first one starts with "110" signaling that the following 8-bit block also belongs to the first and needs to start with "10".

### Discussion

- Which of the following bit strings contain only valid UTF-8 encoded data?
  - 11000000 0100000

No, since the first 8-bit string starts with "110", signaling that the next block also belongs to it, but the next block does not start with "10".

### Discussion

- Which of the following bit strings contain only valid UTF-8 encoded data?
  - 10101010 01011001

No, since the first block starts with "10".

### Discussion

- Which of the following bit strings contain only valid UTF-8 encoded data?
  - 0111

No, only 4 bits.



### Discussion

- Which of the following bit strings contain only valid UTF-8 encoded data?
  - 01001010 01000001

Yes, two 8-bit strings, both starting with 0.

### Discussion

- Which of the following bit strings contain only valid UTF-8 encoded data?
  - 11100010 10011101 10100100

Yes, first string starts with "1110", indicating the next two blocks belong to it. The next two blocks both start with "10".

Questions?

### *B*+-*Trees*

• Issues/Limitations with hash tables?

Range Queries, Space Requirements





### B+-Trees

Solution: Binary Search Trees Even better: B+-trees

• Taking advantage of block accesses

Rules:

- All leaves at same depth
- Terms / postings lists are only at the leaves
- If #keys between d and 2d, then #children between d + 1 and 2d + 1
- Root can have less children

# 1-1 B+-tree 2-3 B+-tree 3-5 B+-tree 4-7 B+-tree

*3-5B+-Tree* 



*3-5B+-Tree* 



*3-5B+-Tree* 





*3-5B+-Tree* 



21.03.2025

8

### *3-5B+-Tree*





*3-5B+-Tree* 





*3-5B+-Tree* 





*3-5B+-Tree* 



*3-5B+-Tree* 



*3-5B+-Tree* 



*3-5B+-Tree* 



*3-5B+-Tree* 



*3-5B+-Tree* 





## Wildcard queries

- e.g. "Z\*rich"
  - Allows different spellings, families of words, etc.
- Leading vs. Trailing queries ("\*ics" vs. "Math\*")
  - "Math\*", just find all matches for Math in the B+tree
  - "\*ics", reverse B+-tree, then do the same

Naïve approach

- Large Storage usage
- > Inefficient
- > False positives



### Permuterm index

### **CHOOSE YOUR RAPPER NAME:**

### Permuterm index





### Permuterm index



k-gram index

1-grams	\$, c, o, m, p, u, t, e, r, \$	Not very
2-grams	\$c, co, om, mp, pu, ut, te,	er, r\$
3-grams	\$co, com, omp, mpu, put,	ute, ter, e
4-grams	\$com, comp, ompu, mput	, pute, ute
5-grams	\$comp, compu, omput, m	pute, pute
6-grams	\$compu, comput, ompute	, mputer, p
7-grams	\$comput, compute, omput	ter, mpute
	1-grams 2-grams 3-grams 4-grams 5-grams 6-grams 7-grams	<ul> <li>1-grams \$, c, o, m, p, u, t, e, r, \$</li> <li>2-grams \$c, co, om, mp, pu, ut, te,</li> <li>3-grams \$co, com, omp, mpu, put,</li> <li>4-grams \$com, comp, ompu, mput</li> <li>5-grams \$comp, compu, omput, m</li> <li>6-grams \$compu, comput, ompute</li> <li>7-grams \$comput, compute, omput</li> </ul>

### useful er\$ er, ter\$ er, uter\$ Not space puter\$ efficient er\$

k-gram index



## Edit/Levenshtein distance (DP...)





## Edit distance (DP...)

- Issue: Costly computations if done on a lot of terms
- Solution: Preselect a few terms
- How?



## Jaccard coefficient



## Jaccard coefficient



Keep terms within large Jaccard coefficients

Soundex algorithm

Change	То
AEHIOUWY	0
BFPV	1
CGJKQSXZ	2
DT	3
L	4
MN	5
R	6



### C0510306 📥 C5136 C513 C510306 **C**5136 **C**513 **Z62** Z06020 Z620 1510650305 1516 **L**36 L0306 L360 R0360104 🗭 R3614 R361

### **Exercise 3: Tolerant Retrieval**

### Bonus Quiz

- Moodle-based
- Start: 21.03 at 15:00
- Deadline: 28.03 at 15:00
- 8 / 16 to pass
- Last two questions involve programming exercises
- The rest are theoretical questions
  - Edit distance
  - Jaccard coefficient
  - B+-trees
- Two attempts in total

- Implement k-gram index
- Process wildcard queries (Bonus)
- Process spell correction (Bonus)

## Wildcard queries

- Implement wildcard\_parse(q) •
- Takes wildcard query ("som\*e") and returns list containing all keys to query the k-gram index



## Wildcard queries

- Implement kgram\_wildcard\_query(q)
- Queries k-gram for matching words, given wildcard query
- Returns list of possible words

```
def kgram_wildcard_query(q):
1
       111
2
       Query the k-gram index for words matching q. Return the matches as a set
3
       111
4
5
       grams = wildcard parse(q)
       ### TODO for the assignment: execute the wildcard query on the k-gram index
6
7
       kgram matches = {}
       post_filter = re.compile(" + q.replace("*", "\\w*") + "$")
8
9
       res = {r for r in kgram matches if post filter.match(r) is not None}
       return res
10
```

## Spell checking

- Modify Jaccard threshold in first code snippet
- Implement spellcorrect(word)
- Compute the k-grams from the input word
- Find matching words for each k-gram
- For each candidate word, check Jaccard coefficient and filter

# https://create.kahoot.it/details/384 <u>02941-280b-43f9-94ba-</u> *4ea8805b9b16*



## FS19 – Question 28

Compute the Jaccard coefficient for the following pairs of words, using their sets of 3-grams (with duplicates eliminated, if any). For example, the set of 3-grams contained in the word "time" is

{"\$ ti", "tim", "ime", "me \$"}

Each answer must be given as an irreducible fraction (or integer).

Question: electric, electronic

Answer:

g(electric) = \$el, ele, lec, ect, ctr, tri, ric, ic\$ g(electronic) = \$el, ele, lec, ect, ctr, tro, ron, oni, nic, ic\$ Intersection: \$el, ele, lec, ect, ctr, ic\$ J(electric, electronic) = 6 / (8 + 10 - 6) = 0.5

## FS19 – Question 28

Compute the Jaccard coefficient for the following pairs of words, using their sets of 3-grams (with duplicates eliminated, if any). For example, the set of 3-grams contained in the word "time" is

```
{"$ ti", "tim", "ime", "me $"}
```

Each answer must be given as an irreducible fraction (or integer) .

```
Question: automobile, automation
```

Answer:

g(automobile) = \$au, aut, uto, tom, omo, mob, obi, bil, ile, le\$ g(automation) = \$au, aut, uto, tom, oma, mat, ati, tio, ion, on\$ Intersection: \$au, aut, uto, tom J(automobile, automation) = 4 / (10 + 10 - 4) = 0.25

## FS19 – Question 28

Compute the Jaccard coefficient for the following pairs of words, using their sets of 3-grams (with duplicates eliminated, if any). For example, the set of 3-grams contained in the word "time" is

{	"\$ti",	"tim",	"ime",	"me\$"	}
---	---------	--------	--------	--------	---

Each answer must be given as an irreducible fraction (or integer).

Question: information, retrieval

Answer:
---------

0

## FS19 – Question 28

Compute the Jaccard coefficient for the following pairs of words, using their sets of 3-grams (with duplicates eliminated, if any). For example, the set of 3-grams contained in the word "time" is

{ "\$ti", "tim", "ime", "me\$" }

Each answer must be given as an irreducible fraction (or integer).

Question: fifteen, fifteen

Answer:

1

### FS19 – Question 32

How many rotations does the following single term contribute to a permuterm index?

mama

\$mama, a\$mam, ma\$ma, ama\$m, mama\$

## FS21 – Question 21

For each one of the following terms, state whether or not it falsely matches the wildcard query «some\*e» if the evaluation of the query is done by simply using a conjunction of 3-grams, and with no post-processing.

sometime

No, since it correctly matches the wildcard query.

## FS21 – Question 21

For each one of the following terms, state whether or not it falsely matches the wildcard query «some\*e» if the evaluation of the query is done by simply using a conjunction of 3-grams, and with no post-processing.

some

Yes, since it is a false positive. (\* is missing)

## FS21 – Question 21

For each one of the following terms, state whether or not it falsely matches the wildcard query «some\*e» if the evaluation of the query is done by simply using a conjunction of 3-grams, and with no post-processing.

sommelier

No, since the 3-gram «ome» doesn't occur in the term it will not be returned.

## FS21 – Question 21

For each one of the following terms, state whether or not it falsely matches the wildcard query «some\*e» if the evaluation of the query is done by simply using a conjunction of 3-grams, and with no post-processing.

something

Yes. «something» covers all 3-grams of the wildcard query, so it is a false positive.